

☐☐ Muon Optimizer: Un Enfoque Basado en Newton-Schulz

☐☐ Integrantes:

- ☐☐ Andrea Parra
- ☐☐ Jorge Andrey

☐☐ Material de apoyo:

- ☐☐ **Diapositivas:** [Ver presentaciones](#)
- ☐☐ **Paper:** [Modular Duality in Deep Learning](#)
- ☐☐ **Código externo:** [TBA](#)
- ☐☐ **Blog muon:** [Muon: An optimizer for hidden layers in neural networks](#)

El optimizador **Muon** es una variante del método de descenso de gradiente con momentum (SGD+Momentum) que introduce una corrección basada en la iteración de **Newton-Schulz** para mejorar la actualización de parámetros. Esta técnica permite estabilizar los gradientes y normalizar la actualización sin necesidad de cálculos costosos como la descomposición en valores singulares (SVD).

En esta sección exploraremos la motivación detrás de Muon, su estructura matemática y cómo Newton-Schulz contribuye a su eficacia.

☐☐ Antecedentes

Este método se basa en el método Newton-Schulz, utilizado para aproximar de la operación Q^{-1} , donde Q es una matriz ortogonal. Sin embargo, los autores del paper proponen una modificación para obtener Q^{-1} , lo que permite la normalización de gradientes en redes neuronales profundas.

Muon busca **aproximar la función de signo de la matriz de gradientes**, lo cual puede interpretarse como "ajustar los valores singulares a 1", asegurando que la actualización de

parámetros mantenga una estructura ortogonal. Este proceso se conoce como "symmetric orthogonalization" y se diferencia del Gram-Schmidt porque no favorece una fila o columna específica.

Los autores eligieron Newton-Schulz por su capacidad para ejecutarse de manera estable en bfloat16, a diferencia de la SVD y otras iteraciones de Newton más costosas o inestables en GPUs modernas.

📄 Objetivos

- ☐ Entender la importancia de métodos de optimización como Muon.
 - ☐ (Intentar) aprender la matemática detrás de la optimización.
 - ☐ Abrir la discusión sobre aplicaciones futuras.
-

📄 Motivación y Contexto

Los métodos de optimización convencionales como SGD y Adam pueden sufrir de **desaparición o explosión del gradiente**, especialmente en redes neuronales profundas. Esto se debe a que la propagación del gradiente puede amplificar valores en ciertas direcciones, afectando la convergencia y estabilidad del entrenamiento.

Los autores de Muon encontraron que en modelos **transformer-based**, las actualizaciones de SGD-momentum y Adam tienen un número de condición muy alto. Es decir, las actualizaciones de pesos están **dominadas por pocas direcciones**, lo que limita la capacidad de aprendizaje. La ortogonalización mediante Newton-Schulz **aumenta la escala de direcciones poco representadas**, ayudando a mejorar la optimización.

⚙️ Método Newton-Schulz en Muon

La iteración de Newton-Schulz es un método iterativo que aproxima la inversa de una matriz sin requerir una factorización directa. En el contexto del optimizador Muon, se utiliza para normalizar el gradiente acumulado BT antes de actualizar los parámetros.

Iteración Newton-Schulz Rectangular

Dado un gradiente acumulado (B_t), se inicia con:

$$X_0 = \frac{B_t}{\|B_t\|_F}$$

Luego, se aplica la siguiente iteración para aproximar UV^T de la SVD de B_t :

$$X_{t+1} = \frac{3}{2} X_t - \frac{1}{2} X_t X_t^T X_t$$

Este método garantiza que los valores singulares se ajusten de manera controlada y que la iteración converja de manera estable, incluso en matrices de bajo rango.

Luego, los autores cambian la fórmula para utilizar una expansión de un polinomio, de la siguiente manera.

$$X_{t+1} = a \cdot X_t + b \cdot X_t X_t^T X_t + c \cdot (X_t X_t^T X_t)^2 X_t + \dots + z \cdot (X_t X_t^T X_t)^n X_t$$

Tomando solo hasta el término cuadrático, y reemplazando X por su SVD, tenemos (en la imagen X es G)

$$G' := aG + b(GG^T)G + c(GG^T)^2 G = \left(aI + b(GG^T) + c(GG^T)^2 \right) G$$

$$= \left(aI + bUS^2U^T + cUS^4U^T \right) USV^T = U \left(aS + bS^3 + cS^5 \right) V^T$$

Luego, si hacemos que la función interna tienda a la función signo, tenemos.

$$G(s) = aS + bS^3 + cS^5 \cdot U \left(\operatorname{sig}(G(s)) \right) V^T$$

Si la función signo tiende a 1, tenemos una aproximación a UV^T , que es la matriz ortogonal del gradiente.

Los autores también exploraron el ajuste de coeficientes de Newton-Schulz para acelerar la convergencia, logrando reducir la cantidad de iteraciones necesarias a solo 5 en sus experimentos.

□□ Diferencias entre Muon y SGD+Momentum

La principal diferencia de Muon con SGD+Momentum es la inclusión del término (\mathbf{o}_t), obtenido mediante Newton-Schulz:

SGD+Momentum (Convencional)

```
Require: Learning rate  $\eta$ , momentum  $\mu$ 
Inicializar parámetros  $\theta_0$ 
Inicializar velocidad  $v_0 \leftarrow 0$ 
for  $t = 1, \dots$  do
  Calcular gradiente  $G_t \leftarrow \nabla \theta \mathcal{L}_t(\theta_{t-1})$ 
  Actualizar velocidad  $v_t \leftarrow \mu v_{t-1} + G_t$ 
  Actualizar parámetros  $\theta_t \leftarrow \theta_{t-1} - \eta v_t$ 
end for
return  $\theta_t$ 
```

Muon Optimizer

```
Require: Learning rate  $\eta$ , momentum  $\mu$ 
Inicializar  $B_0 \leftarrow 0$ 
for  $t = 1, \dots$  do
  Calcular gradiente  $G_t \leftarrow \nabla \theta \mathcal{L}_t(\theta_{t-1})$ 
   $B_t \leftarrow \mu B_{t-1} + G_t$ 
   $O_t \leftarrow \text{NewtonSchulz5}(B_t)$ 
  Actualizar parámetros  $\theta_t \leftarrow \theta_{t-1} - \eta O_t$ 
end for
return  $\theta_t$ 
```

- En **SGD+Momentum**, el gradiente acumulado se usa directamente para actualizar los parámetros.
- En **Muon**, el gradiente acumulado B_t es corregido mediante Newton-Schulz para obtener O_t , asegurando una actualización bien condicionada.

Esto ayuda a evitar direcciones de gradiente mal condicionadas y estabiliza la convergencia.

☐☐ Aplicaciones y Beneficios

El uso de Muon y Newton-Schulz tiene aplicaciones en diversas áreas:

- **Redes neuronales profundas:** Mejora la estabilidad en entrenamientos largos.
- **Redes recurrentes (RNNs, Transformers):** Evita problemas de explosión/desaparición del gradiente.
- **Generative Adversarial Networks (GANs):** Regulariza el entrenamiento para mejorar la calidad de las muestras generadas.
- **Optimización de alto rendimiento:** Reduce la necesidad de ajustes manuales en el learning rate.

□ □ **Ventajas clave de Muon:**

- ✓ Evita el costo computacional de la SVD.
- ✓ Mantiene estabilidad en gradientes.
- ✓ Regulariza la actualización de parámetros.

□ □ **Referencias**

- □ [Webpage autores de Muon sobre Newton-Schulz](#)
- □ [Explicación en twitter](#)
- □ Higham, N. J. Functions of Matrices. SIAM, 2008.

Revision #28

Created 7 March 2025 20:30:42 by Jorge Garcia

Updated 14 August 2025 18:42:42 by Jorge Garcia